

# Parallel System Performance: Evaluation & Scalability

- **Factors affecting parallel system performance:**
  - Algorithm-related, parallel program related, architecture/hardware-related.
- **Workload-Driven Quantitative Architectural Evaluation:**
  - Select applications or suite of parallel benchmarks to evaluate parallel architecture either on real or simulated machine.
  - From measured performance results compute parallel performance metrics:
    - Speedup, System Efficiency, Redundancy, Utilization, Quality of Parallelism.
  - **Resource-oriented Workload scaling models:** How the speedup of a parallel computation is affected subject to specific constraints:
    - 1 • *Problem constrained* (PC): Fixed-load Model.
    - 2 • *Time constrained* (TC): Fixed-time Model.
    - 3 • *Memory constrained* (MC): Fixed-Memory Model. Data usage per processor is fixed
- **Parallel Performance Scalability:**
  - Definition.
  - Conditions of scalability.
  - Factors affecting scalability.

Informally:  
The ability of parallel system performance to increase with increased problem size and system size.

**For a given parallel system and a given parallel computation/problem/algorithm**

# Parallel Program Performance

- Parallel processing goal is to maximize speedup:

$$\text{Speedup} = \frac{\text{Time}(1)}{\text{Time}(p)} \leq \frac{\text{Sequential Work (or Time)}}{\text{Max (Work + Synch Wait Time + Comm Cost + Extra Work)}}$$

Fixed Problem Size Speedup

Max for any processor

Parallelizing Overheads

- By:

- 1 – Balancing computations/overheads (workload) on processors (every processor has the same amount of work/overheads).
- 2 – Minimizing communication cost and other overheads associated with each step of parallel program creation and execution.

## Parallel Performance Scalability:

For a given parallel system and parallel computation/problem/algorithm

Achieve a good speedup for the parallel application on the parallel architecture as problem size and machine size (number of processors) are increased.

Or

Continue to achieve good parallel performance "speedup" as the sizes of the system/problem are increased.

(More formal treatment of scalability later)

# Factors affecting Parallel System Performance

## • Parallel Algorithm-related:

i.e. Max.  
Inherent  
Parallelism

- Available concurrency and profile, dependency graph, uniformity, patterns.
- Complexity and predictability of computational requirements.
- Required communication/synchronization, uniformity and patterns.
- Data size requirements. - Type of parallelism: Data and/or Functional.

## • Parallel program related:

- Partitioning: Decomposition and assignment to tasks
  - Parallel task grain size.
  - Communication to computation ratio. C-to-C ratio (measure of inherent communication)  
For a given partition
- Programming model used.
- Orchestration
  - Cost of communication/synchronization.
- Resulting data/code memory requirements, locality and working set characteristics.
- Mapping & Scheduling: Dynamic or static.

## • Hardware/Architecture related:

- Total CPU computational power available. + Number of Processors
- Parallel programming model support:
  - e.g support for Shared address space Vs. message passing support.
  - Architectural interactions, artifactual “extra” communication
- Communication network characteristics: Scalability, topology ..
- Memory hierarchy properties.

Refined from factors in Lecture # 1

**CMPE655 - Shaaban**

# Parallel Performance Metrics Revisited

Observed Concurrency Profile

MIN( Software Parallelism , Hardware Parallelism )

- **Degree of Parallelism (DOP):** For a given time period, reflects the number of processors in a specific parallel computer actually executing a particular parallel program.

- **Average Parallelism, A:** i.e DOP at a given time = Min (Software Parallelism, Hardware Parallelism)

- Given maximum parallelism =  $m$
- $n$  homogeneous processors
- Computing capacity of a single processor  $\Delta$  ↙ Computations/sec
- Total amount of work (instructions or computations):

$$W = \Delta \int_{t_1}^{t_2} DOP(t) dt \quad \text{or as a discrete summation} \quad W = \Delta \sum_{i=1}^m i \cdot t_i$$

Where  $t_i$  is the total time that  $DOP = i$  and  $\sum_{i=1}^m t_i = t_2 - t_1$  ↙ Execution Time

**The average parallelism A:**

$$A = \frac{1}{t_2 - t_1} \int_{t_1}^{t_2} DOP(t) dt$$

In discrete form

$$A = \frac{\left( \sum_{i=1}^m i \cdot t_i \right)}{\left( \sum_{i=1}^m t_i \right)}$$

Execution Time →

DOP Area

Total Work ↗

↖ Execution Time

From Lecture # 3

**CMPE655 - Shaaban**

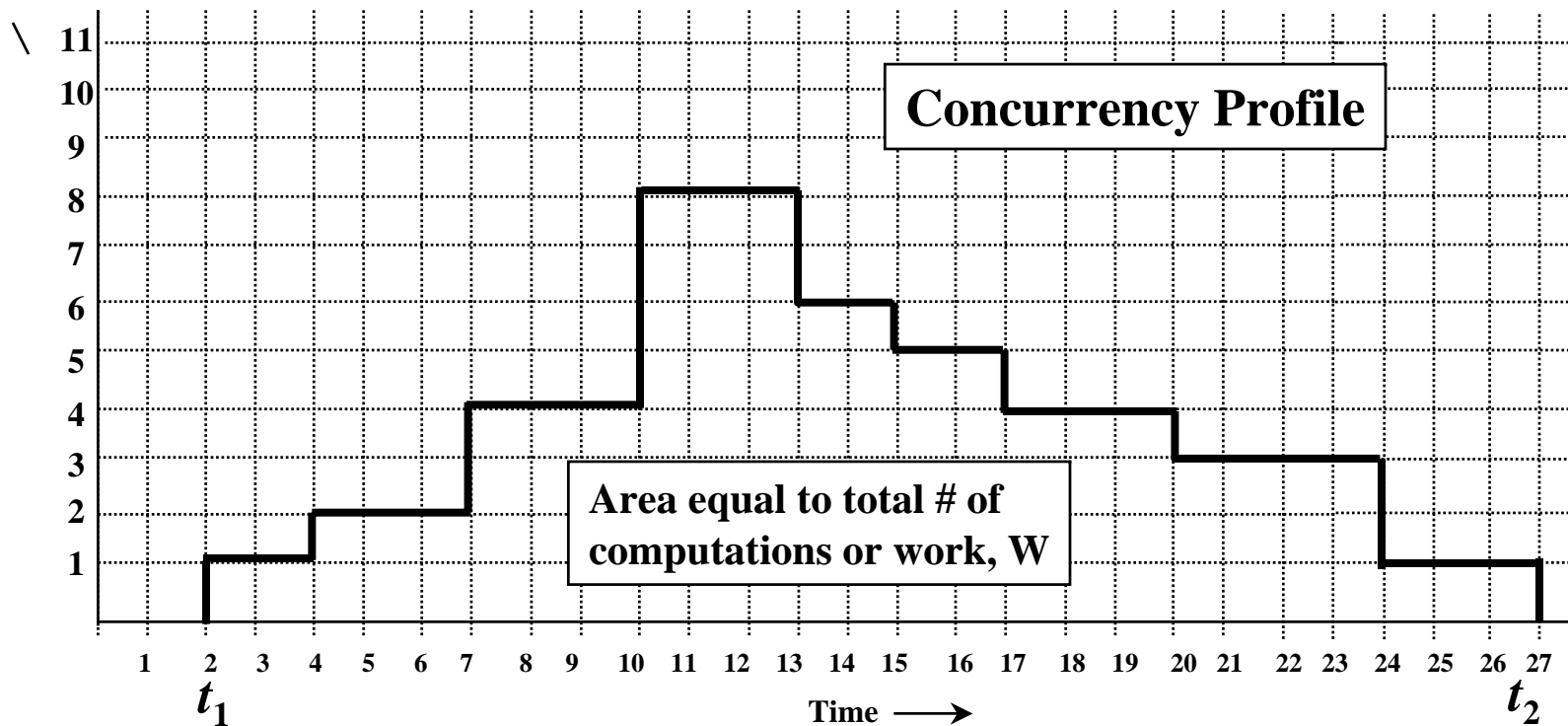
# Example: Concurrency Profile of A Divide-and-Conquer Algorithm

- Execution observed from  $t_1 = 2$  to  $t_2 = 27$
- Peak parallelism  $m = 8$
- $A = (1 \times 5 + 2 \times 3 + 3 \times 4 + 4 \times 6 + 5 \times 2 + 6 \times 2 + 8 \times 3) / (5 + 3 + 4 + 6 + 2 + 2 + 3)$   
 $= 93/25 = 3.72$

$$A = \left( \sum_{i=1}^m i \cdot t_i \right) / \left( \sum_{i=1}^m t_i \right)$$

Average Parallelism

Degree of Parallelism (DOP)



From Lecture # 3

CMPE655 - Shaaban

# Parallel Performance Metrics Revisited

**Asymptotic Speedup:**

i.e. Hardware Parallelism  $n >$  Software Parallelism  $m$

(more processors  $n$  than max software DOP,  $m$ )

**Execution time with one processor**

$$T(1) = \sum_{i=1}^m t_i(1) = \sum_{i=1}^m \frac{W_i}{\Delta}$$

**Execution time with an infinite number of available processors (number of processors  $n = \infty$  or  $n > m$ )**

$$T(\infty) = \sum_{i=1}^m t_i(\infty) = \sum_{i=1}^m \frac{W_i}{i\Delta}$$

Sequential Work or Time

**Asymptotic speedup  $S_\infty$**

$$S_\infty = \frac{T(1)}{T(\infty)} = \frac{\sum_{i=1}^m W_i}{\sum_{i=1}^m W_i / i}$$

**The above ignores all overheads.**

- $\Delta$  = Computing capacity of a single processor
- $m$  = maximum degree of software parallelism
- $t_i$  = total time that DOP =  $i$
- $W_i$  = total work with DOP =  $i$

Keeping problem size fixed and ignoring parallelization overheads/extra work

**CMPE655 - Shaaban**

i.e. Hardware parallelism  $n$  exceeds software parallelism  $m$

# Phase Parallel Model of An Application

- Consider a sequential program of size  $s$  consisting of  $k$  computational phases  $C_1 \dots C_k$  where each phase  $C_i$  has a degree of parallelism  $DOP = i$
- Assume single processor execution time of phase  $C_i = T_1(i)$
- Total single processor execution time =  $T_1 = \sum_{i=1}^{i=k} T_1(i)$  k = max. DOP
- Ignoring overheads,  $n$  processor execution time:  $T_n = \sum_{i=1}^{i=k} T_1(i) / \min(i, n)$
- If all overheads are grouped as interaction  $T_{\text{interact}} = \text{Synch Time} + \text{Comm Cost}$  and parallelism  $T_{\text{par}} = \text{Extra Work}$ , as  $h(s, n) = T_{\text{interact}} + T_{\text{par}}$  then parallel execution time:

Accounting for parallelization overheads

Lump sum overheads term  $h(s, n)$

$$T_n = \sum_{i=1}^{i=k} T_1(i) / \min(i, n) + h(s, n)$$

Total overheads

$n$  = number of processors

$s$  = problem size

DOP Profile

- If  $k = n$  and  $f_i$  is the fraction of sequential execution time with  $DOP = i$   $\pi = \{f_i | i = 1, 2, \dots, n\}$  and ignoring overheads ( $h(s, n) = 0$ ) the speedup is given by:

$$S(n) = S(\infty) = \frac{T_1}{T_n} = \frac{1}{\left(\sum_{i=1}^n f_i / i\right)}$$

$\pi = \{f_i | i = 1, 2, \dots, n\}$  for max DOP =  $n$   
is parallelism degree probability distribution (DOP profile)

CMPE655 - Shaaban

# Harmonic Mean Speedup for $n$ Execution Mode Multiprocessor system

Fig 3.2 page 111  
See handout



# Parallel Performance Metrics Revisited: Amdahl's Law

- **Harmonic Mean Speedup** ( $i$  number of processors used  $f_i$  is the fraction of sequential execution time with DOP =  $i$ ):

$$S(n) = T_1 / T_n = \frac{1}{\left( \sum_{i=1}^n f_i / i \right)}$$

DOP = 1  
(sequential)

DOP =  $n$

- In the case  $\pi = \{f_i \text{ for } i = 1, 2, \dots, n\} = (\alpha, 0, 0, \dots, 1-\alpha)$ , the system is running sequential code with probability  $\alpha$  and utilizing  $n$  processors with probability  $(1-\alpha)$  with other processor modes not utilized (i.e. = 0).

**Amdahl's Law:**

$$S_n = \frac{1}{\alpha + (1-\alpha)/n}$$

Keeping problem size fixed  
and ignoring overheads  
(i.e.  $h(s, n) = 0$ )

+  $h(s, n)$  ?

$S \rightarrow 1/\alpha$  as  $n \rightarrow \infty$

⇒ Under these conditions the best speedup is upper-bounded by  $1/\alpha$

CMPE655 - Shaaban

Alpha =  $\alpha$  = Sequential fraction with DOP = 1

# Parallel Performance Metrics Revisited

## Efficiency, Utilization, Redundancy, Quality of Parallelism

- **System Efficiency:** Let  $O(n)$  be the total number of unit operations performed by an  $n$ -processor system and  $T(n)$  be the parallel execution time in unit time steps:
  - In general  $T(n) \ll O(n)$  (more than one operation is performed by more than one processor in unit time).

— Assume  $T(1) = O(1)$  For One Processor

— **Speedup factor:**  $S(n) = T(1) / T(n)$

- Ideal  $T(n) = T(1)/n \rightarrow$  Ideal speedup =  $n$

— **Parallel System efficiency**  $E(n)$  for an  $n$ -processor system:

$$E(n) = S(n)/n = T(1)/[nT(n)]$$

**Ideally:**

**Ideal speedup:**  $S(n) = n$

**and thus ideal efficiency:**  $E(n) = n/n = 1$

$n$  = number of processors

Here  $O(1)$  = work on one processor

$O(n)$  = total work on  $n$  processors

# Parallel Performance Metrics Revisited

## Cost, Utilization, Redundancy, Quality of Parallelism

- **Cost**: The processor-time product or cost of a computation is defined as:

$$\text{Speedup} = T(1)/T(n)$$

$$\text{Efficiency} = S(n)/n$$

$$\text{Cost}(n) = n T(n) = n \times T(1) / S(n) = T(1) / E(n)$$

- The cost of sequential computation on one processor  $n=1$  is simply  $T(1)$
- A cost-optimal parallel computation on  $n$  processors has a cost proportional to  $T(1)$  when:

Ideal parallel speedup

$$S(n) = n, E(n) = 1 \quad \text{--->} \quad \text{Cost}(n) = T(1)$$

- **Redundancy**:  $R(n) = O(n)/O(1)$

- Ideally with no overheads/extra work  $O(n) = O(1) \rightarrow R(n) = 1$

- **Utilization**:  $U(n) = R(n)E(n) = O(n) / [nT(n)]$

- ideally  $R(n) = E(n) = U(n) = 1$

Perfect load balance?

Assuming:  
 $T(1) = O(1)$

- **Quality of Parallelism**:

$$Q(n) = S(n) E(n) / R(n) = T^3(1) / [nT^2(n)O(n)]$$

- Ideally  $S(n) = n, E(n) = R(n) = 1 \quad \text{--->} \quad Q(n) = n$

$n$  = number of processors

here:  $O(1)$  = work on one processor  $O(n)$  = total work on  $n$  processors

CMPE655 - Shaaban

# A Parallel Performance measures Example

For a hypothetical workload with

- $O(1) = T(1) = n^3$  Work or time on one processor
- $O(n) = n^3 + n^2 \log_2 n$       $T(n) = 4n^3 / (n+3)$   
Total parallel work on n processors     Parallel execution time on n processors
- $\text{Cost}(n) = 4n^4 / (n+3) \sim 4n^3$

**Fig 3.4 page 114**

---

**Table 3.1 page 115**  
**See handout**

# Application Scaling Models for Parallel Computing

- If work load  $W$  or problem size “ $s$ ” is unchanged then: **i.e. Fixed problem size**
  - The efficiency  $E$  may decrease as the machine size  $n$  increases if the overhead  $h(s, n)$  increases faster than the increase in machine size.
- The condition of a scalable parallel computer solving a scalable parallel problem exists when:
  - A desired level of efficiency is maintained by increasing the machine size “ $n$ ” and problem size “ $s$ ” proportionally.  **$E(n) = S(n)/n$**
  - In the ideal case the workload curve is a linear function of  $n$ : (*Linear scalability in problem size*).
- Application Workload Scaling Models for Parallel Computing:

Workload scales subject to a given constraint as the machine size is increased:

- 1 – Problem constrained (PC): or Fixed-load Model. Corresponds to a constant workload or fixed problem size.
- 2 – Time constrained (TC): or Fixed-time Model. Constant execution time.
- 3 – Memory constrained (MC): or Fixed-memory Model: Scale problem so memory usage per processor stays fixed. Bound by memory of a single processor.

Data usage per processor is fixed

**What about Iso-Efficiency? (Fixed Efficiency?) ....**

**CMPE655 - Shaaban**

**$n$  = Number of processors       $s$  = Problem size**

# Problem Constrained (PC) Scaling : Fixed-Workload Speedup

When  $DOP = i > n$  ( $n = \text{number of processors}$ )

$i = 1 \dots m$  i.e.  $n > m$

Execution time of  $W_i$

$$t_i(n) = \frac{W_i}{i\Delta} \left\lceil \frac{i}{n} \right\rceil$$

Total execution time

$$T(n) = \sum_{i=1}^m \frac{W_i}{i\Delta} \left\lceil \frac{i}{n} \right\rceil$$

Corresponds to "Normal" parallel speedup:  
Keep problem size (workload) fixed as the size of the parallel machine (number of processors) is increased.

If  $DOP = i < n$ , then  $t_i(n) = t_i(\infty) = W_i/i\Delta$

Ignoring parallelization overheads

$S_n$

Fixed-load speedup factor is defined as the ratio of  $T(1)$  to  $T(n)$ :

Ignoring overheads

$$S_n = \frac{T(1)}{T(n)} = \frac{\sum_{i=1}^m W_i}{\sum_{i=1}^m \frac{W_i}{i} \left\lceil \frac{i}{n} \right\rceil}$$

Let  $h(s, n)$  be the total system overheads on an  $n$ -processor system:

$$S_n = \frac{T(1)}{T(n) + h(s, n)} = \frac{\sum_{i=1}^m W_i}{\sum_{i=1}^m \frac{W_i}{i} \left\lceil \frac{i}{n} \right\rceil + h(s, n)}$$

The overhead term  $h(s, n)$  is both application- and machine-dependent and usually difficult to obtain in closed form.

Total parallelization overheads term

# Amdahl's Law for Fixed-Load Speedup

- For the special case where the system either operates in sequential mode (DOP = 1) or a perfect parallel mode (DOP =  $n$ ), the Fixed-load speedup is simplified to:

$n$  = number of processors

$$S_n = \frac{W_1 + W_n}{W_1 + W_n/n}$$

i.e. ignoring parallelization overheads

We assume here that the overhead factor  $h(s, n) = 0$

For the normalized case where: / Sequential Fraction DOP = 1

$$W_1 + W_n = \alpha + (1 - \alpha) = 1 \text{ with } \alpha = W_1 \text{ and } 1 - \alpha = W_n \setminus$$

Parallel Fraction DOP =  $n$

The equation is reduced to the previously seen form of

Amdahl's Law:

$$S_n = \frac{1}{\alpha + (1 - \alpha)/n}$$

CMPE655 - Shaaban

Alpha =  $\alpha$  = Sequential fraction with DOP = 1

# Time Constrained (TC) Workload Scaling

Both problem size (workload) and machine size are scaled (increased) so execution time remains constant.

## $S'_n$ Fixed-Time Speedup

- To run the largest problem size possible on a larger machine with about the same execution time of the original problem on a single processor.

Let  $m'$  be the maximum DOP for the scaled up problem,

$W'_i$  be the scaled workload with DOP =  $i$

In general,  $W'_i > W_i$  for  $2 \leq i \leq m'$  and  $W'_1 = W_1$  assumption

Assuming that  $T(1) = T'(n)$  we obtain :

**i.e fixed execution time**  $T(1) = T'(n) = \sum_{i=1}^m W_i = \sum_{i=1}^{m'} \frac{W'_i}{i} \left\lceil \frac{i}{n} \right\rceil + h(s, n)$

$S'_n$  Speedup  $S'_n = T'(1) / T'(n)$  is given by:

$$S'_n = \frac{T'(1)}{T'(n)} = \frac{T'(1)}{T(1)} = \frac{\sum_{i=1}^{m'} W'_i}{\sum_{i=1}^m \frac{W'_i}{i} \left\lceil \frac{i}{n} \right\rceil + h(s, n)} = \frac{\sum_{i=1}^{m'} W'_i}{\sum_{i=1}^m W_i}$$

Time on one processor for scaled problem

Original workload

**Fixed-Time Speedup**

Total parallelization overheads term

**CMPE655 - Shaaban**

$s$  = problem size  $n$  = number of processors



# Gustafson's Fixed-Time Speedup

- For the special fixed-time speedup case where DOP can either be 1 or  $n$  and assuming  $h(s,n) = 0$  i.e no overheads

Time for scaled up problem on one processor fixed execution time  $T(1) = T'(n)$  Also assuming:  $W'_1 = W_1$

$$S'_n = \frac{T'(1)}{T'(n)} = \frac{T'(1)}{T(1)} = \frac{\sum_{i=1}^{m'} W'_i}{\sum_{i=1}^m W_i} = \frac{W'_1 + W'_n}{W_1 + W_n} = \frac{W_1 + nW_n}{W_1 + W_n}$$

assumption

For Original Problem Size

Where  $W'_n = nW_n$  and  $T(1) = W_1 + W_n = W'_1 + W'_n/n = T'(n)$

DOP = 1

DOP = n

Assuming  $\alpha = W_1$  and  $1 - \alpha = W_n$  and  $W_1 + W_n = 1$

(i.e normalize to 1)

$S'_n$

$$S'_n = \frac{T'(1)}{T'(n)} = \frac{\alpha + n(1 - \alpha)}{\alpha + (1 - \alpha)} = n - \alpha(n - 1)$$

$T(1) = T'(n)$

CMPE655 - Shaaban

Alpha =  $\alpha$  = Sequential fraction with DOP = 1

# Memory Constrained (MC) Scaling

## Fixed-Memory Speedup

Problem and machine size

Problem and machine sizes are scaled so memory usage per processor  $M$  stays fixed.

- Scale so memory usage per processor stays fixed
- **Scaled Speedup:  $\text{Time}(1) / \text{Time}(n)$**  for scaled up problem
- Let  $M$  be the memory requirement of a given problem
- Let  $W = g(M)$  or  $M = g^{-1}(W)$  where

**Scaled up problem memory requirement =  $nM$**   
 **$n$  = number of processors**  
 **$M$  = memory requirement for one processor**

$$W = \sum_{i=1}^m W_i \text{ workload for sequential execution} \quad W^* = \sum_{i=1}^m W_i^* \text{ scaled workload on } n \text{ nodes}$$

The memory bound for an active node is  $g^{-1}\left(\sum_{i=1}^m W_i\right)$

**The fixed-memory speedup is defined by:**

Assuming  $W_n^* = g^{-1}(nM) = G(n)g(M) = G(n)W_n$   
 and either sequential or perfect parallelim and  $h(s, n) = 0$

$$S_n^* = \frac{T^*(1)}{T^*(n)} = \frac{\sum_{i=1}^m W_i^*}{\sum_{i=1}^m \frac{W_i^*}{i} \left\lceil \frac{i}{n} \right\rceil + h(s, n)}$$

**DOP = 1**      **DOP = n**      **No overheads**

$$S_n^* = \frac{W_1^* + W_n^*}{W_1^* + W_n^*/n} = \frac{W_1 + G(n)W_n}{W_1 + G(n)W_n/n}$$

**4 cases for  $G(n)$**       **Also assuming:  $W_1^* = W_1$**

- 1
- 2
- 3
- 4

- $G(n) = 1$  problem size fixed (Amdahl's)**
- $G(n) = n$  workload increases  $n$  times as memory demands increase  $n$  times = Fixed Time**
- $G(n) > n$  workload increases faster than memory requirements  $S_n^* > S'_n$  — Fixed-Time Speedup**
- $G(n) < n$  memory requirements increase faster than workload  $S'_n > S_n^*$  — Fixed-Memory Speedup**

$S_n^*$  Memory Constrained, MC (fixed memory) speedup  
 $S'_n$  Time Constrained, TC (fixed time) speedup

**CMPE655 - Shaaban**

# Impact of Scaling Models: 2D Grid Solver

- For sequential  $n \times n$  solver: memory requirements  $O(n^2)$ . Computational complexity  $O(n^2)$  times number of iterations (minimum  $O(n)$ ) thus  $W = O(n^3)$

1

## Problem constrained (PC) Scaling:

Fixed problem size

Number of iterations

Total work

– Grid size fixed =  $n \times n$       Ideal Parallel Execution time =  $O(n^3/p)$

Parallelization overheads ignored

– Memory requirements per processor =  $O(n^2/p)$

2

## Memory Constrained (MC) Scaling:

– Memory requirements stay the same:  $O(n^2)$  per processor.

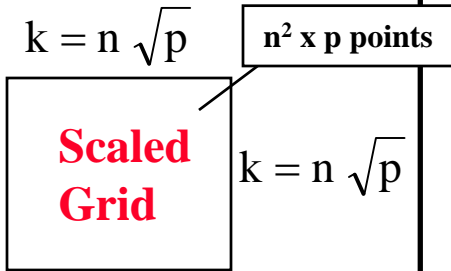
– Scaled grid size =  $k \times k = n \sqrt{p} \times n \sqrt{p} = n^2 p$  points

– Iterations to converge =  $n \sqrt{p} = k$  (new grid size)

– Workload =  $O\left((n \sqrt{p})^3\right)$

– Ideal parallel execution time =  $O\left(\frac{(n \sqrt{p})^3}{p}\right) = O\left(n^3 \sqrt{p}\right)$

• Grows by  $\sqrt{p}$



Parallelization overheads ignored

**Example:**

1 hr on uniprocessor for original problem means 32 hr on 1024 processors for scaled up problem (new grid size  $32n \times 32n$ ).

3

## Time Constrained (TC) scaling:

– Execution time remains the same  $O(n^3)$  as sequential case.

– If scaled grid size is  $k \times k$ , then  $k^3/p = n^3$ , so  $k = n \times \sqrt[3]{p}$

– Memory requirements per processor =  $k^2/p = \frac{n^2}{\sqrt[3]{p}}$

• Diminishes as cube root of number of processors

Workload =

$$O\left((n \sqrt[3]{p})^3\right) = O(n^3 p)$$

Grows slower than MC

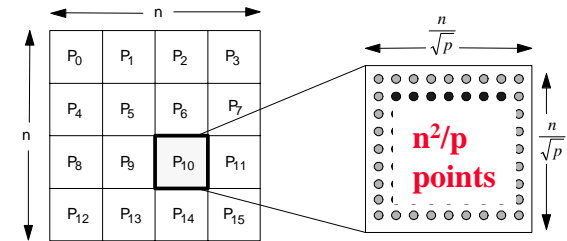
**CMPE655 - Shaaban**

$p$  = number of processors       $n \times n$  = original grid size

# Impact on Grid Solver Execution Characteristics

- **Maximum Concurrency: Total Number of Grid points**

- PC: fixed;  $n^2$
- MC: grows as  $p$ :  $p \times n^2$
- TC: grows as  $p^{0.67}$   $n^2 \times p^{\frac{2}{3}}$



- **Comm. to comp. Ratio: Assuming block decomposition**

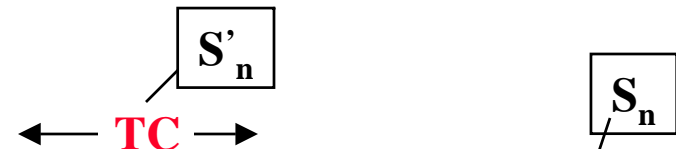
- PC: grows as  $\sqrt{p}$ ; Grid size  $n$  fixed     $Communication = \frac{4n}{\sqrt{p}}$      $Computation = \frac{n^2}{p}$
- MC: fixed;  $4/n$     New grid size  $k = n \sqrt{p}$
- TC: grows as  $\sqrt[6]{p}$     New grid size  $k = n \times \sqrt[3]{p}$      $original \ c-to-c = \frac{4 \times \sqrt{p}}{n}$

- **Working Set: (i.e. Memory requirements per processor)**

PC: shrinks as  $p$  :  $n^2/p$

MC: fixed =  $n^2$

TC: shrinks as  $\sqrt[3]{p}$  :  $n^2/\sqrt[3]{p}$



- **Expect speedups to be best under MC and worst under PC.**

PC= Problem constrained = Fixed-load or fixed problem size model  
 MC = Memory constrained = Fixed-memory Model  
 TC = Time constrained = Fixed-time Model

S<sub>n</sub>\*

CMPE655 - Shaaban

For a given parallel system and parallel computation/problem/algorithm

# Scalability

... of Parallel Architecture/Algorithm Combination

- The study of scalability in parallel processing is concerned with determining the degree of matching between a parallel computer architecture and application/algorithm and whether this degree of matching continues to hold as problem and machine sizes are scaled up.
- Combined architecture/algorithmic scalability imply increased problem size can be processed with acceptable performance level with increased system size for a particular architecture and algorithm.

**i.e.** → – Continue to achieve good parallel performance "speedup" as the sizes of the system/problem are increased.

- Basic factors affecting the scalability of a parallel system for a given problem:

Machine Size  $n$

Clock rate  $f$

Problem Size  $s$

CPU time  $T$

I/O Demand  $d$

Memory Capacity  $m$

Communication/other overheads  $h(s, n)$ , where  $h(s, 1) = 0$

Computer Cost  $c$

Programming Overhead  $p$

For scalability, overhead term must grow slowly as problem/system sizes are increased

Parallel Architecture ↔ Parallel Algorithm  
Match?

As sizes increase

CMPE655 - Shaaban

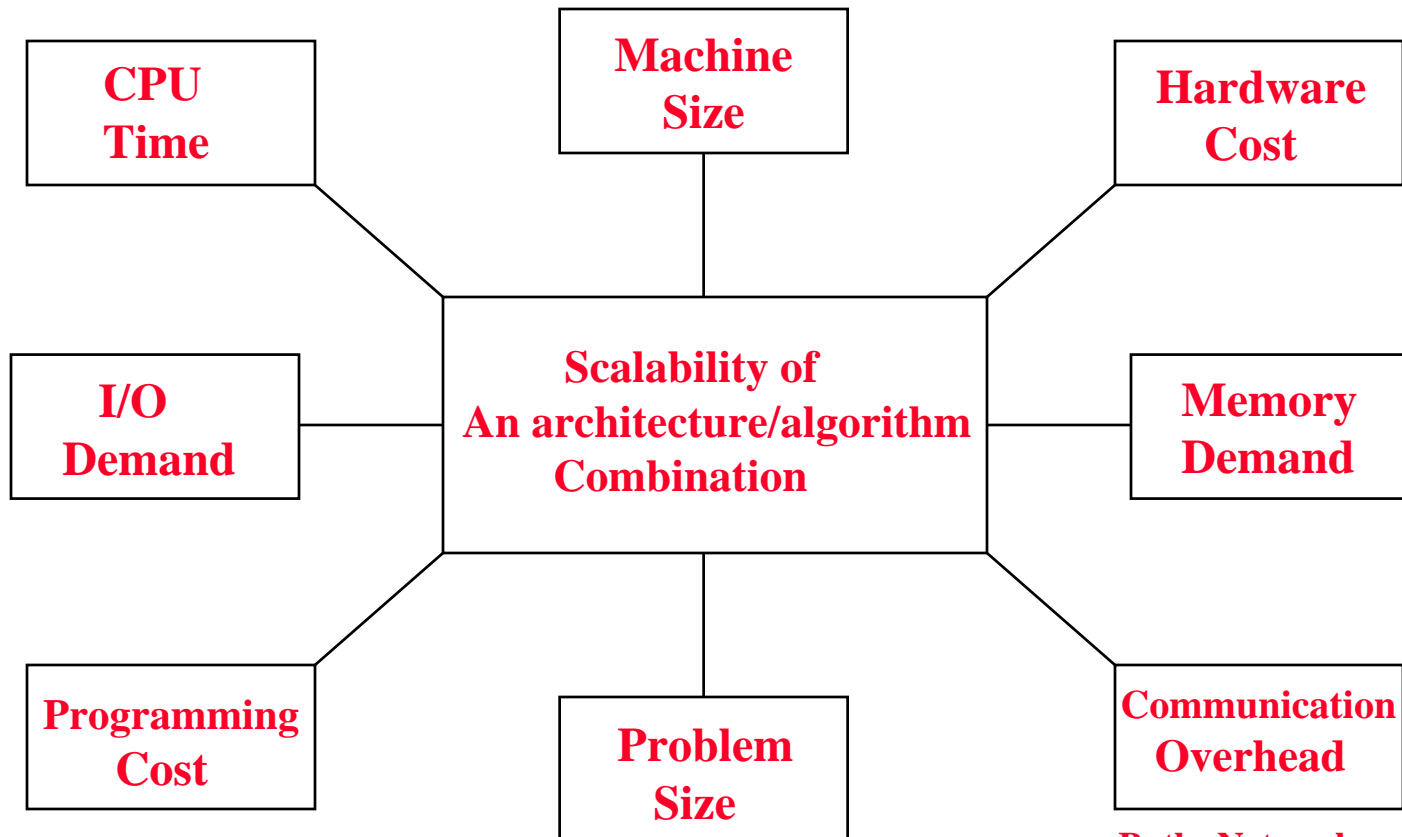
# Parallel Scalability Factors

•The study of scalability in parallel processing is concerned with determining the degree of matching between a parallel computer architecture and application/algorithm and whether this degree of matching continues to hold as problem and machine sizes are scaled up .

•Combined architecture/algorithmic scalability imply increased problem size can be processed with acceptable performance level with increased system size for a particular architecture and algorithm.

–Continue to achieve good parallel performance "speedup" as the sizes of the system/problem are increased.

From last slide



Both: Network + software overheads

For a given parallel system and parallel computation/problem/algorithm

**CMPE655 - Shaaban**

# Revised Asymptotic Speedup, Efficiency

Vary both problem size  $S$  and number of processors  $n$

## • Revised Asymptotic Speedup:

Accounting for overheads

$$S(s, n) = \frac{T(s, 1)}{T(s, n) + h(s, n)}$$

Condition for scalability

Problem/Architecture  
Scalable if  $h(s, n)$  grows  
slowly as  $s, n$  increase

Based on DOP profile

- $s$  problem size.
- $n$  number of processors
- $T(s, 1)$  minimal sequential execution time on a uniprocessor.
- $T(s, n)$  minimal parallel execution time on an  $n$ -processor system.
- $h(s, n)$  lump sum of all communication and other overheads.

## • Revised Asymptotic Efficiency:

$$E(s, n) = \frac{S(s, n)}{n}$$

Iso-Efficiency? (Fixed Efficiency?)

CMPE655 - Shaaban

# Parallel System Scalability

- Scalability (very restrictive definition):

A system architecture is scalable if the system efficiency  $E(s, n) = 1$  for all algorithms with any number of processors  $n$  and any size problem  $s$

- Another Scalability Definition (more formal, less restrictive):

The scalability  $\Phi(s, n)$  of a machine for a given algorithm is defined as the ratio of the asymptotic speedup  $S(s, n)$  on the real machine to the asymptotic speedup  $S_I(s, n)$

“Ideal” PRAM Speedup

$$S_I(s, n) = \frac{T(s, 1)}{T_I(s, n)}$$

on the ideal realization of an EREW PRAM

For real parallel machine

For PRAM

$$\Phi(s, n) = \frac{S(s, n)}{S_I(s, n)} = \frac{T_I(s, n)}{T(s, n)}$$

Capital Phi

For PRAM

For real parallel machine

Ideal  $\Phi$  ?

CMPE655 - Shaaban

$s$  = size of problem     $n$  = number of processors



# **Example: Scalability of Network Architectures for Parity Calculation**

**Table 3.7 page 142  
see handout**

# Evaluating a Real Parallel Machine

- Performance Isolation using Microbenchmarks
- Choosing Workloads
- Evaluating a Fixed-size Machine
- Varying Machine Size and Problem Size
- All these issues, plus more, relevant to evaluating a tradeoff via simulation

**To Evaluate  
Scalability**

# Performance Isolation: Microbenchmarks

- **Microbenchmarks**: Small, specially written programs to isolate performance characteristics
  - Processing.
  - Local memory.
  - Input/output.
  - Communication and remote access (read/write, send/receive).
  - Synchronization (locks, barriers).
  - Contention.
  - Network
  - .....

# Types of Workloads/Benchmarks

- *Kernels*: matrix factorization, FFT, depth-first tree search
- *Complete Applications*: ocean simulation, ray trace, database.
- *Multiprogrammed Workloads*.

• Multiprog. ↔ Appls ↔ Kernels ↔ Microbench.

Realistic  
Complex  
Higher level interactions  
Are what really matters

Easier to understand  
Controlled  
Repeatable  
Basic machine characteristics

Each has its place:

*Use kernels and microbenchmarks to gain understanding, but full applications needed to evaluate realistic effectiveness and performance*

# **Three Desirable Properties for Parallel Workloads**

- 1. Representative of application domains.**
- 2. Coverage of behavioral properties.**
- 3. Adequate concurrency.**

## Desirable Properties of Workloads:

### 1 Representative of Application Domains

- Should adequately represent domains of interest, e.g.:
  - *Scientific*: Physics, Chemistry, Biology, Weather ...
  - *Engineering*: CAD, Circuit Analysis ...
  - *Graphics*: Rendering, radiosity ...
  - *Information management*: Databases, transaction processing, decision support ...
  - *Optimization*
  - *Artificial Intelligence*: Robotics, expert systems ...
  - *Multiprogrammed general-purpose workloads*
  - *System software*: e.g. the operating system

Etc....

CMPE655 - Shaaban

# Desirable Properties of Workloads:

## 2 Coverage: Stressing Features

- **Some features of interest to be covered by workload:**
  - **Compute v. memory v. communication v. I/O bound**
  - **Working set size and spatial locality**
  - **Local memory and communication bandwidth needs**
  - **Importance of communication latency**
  - **Fine-grained or coarse-grained**
    - **Data access, communication, task size**
  - **Synchronization patterns and granularity**
  - **Contention**
  - **Communication patterns**

- **Choose workloads that cover a range of properties**

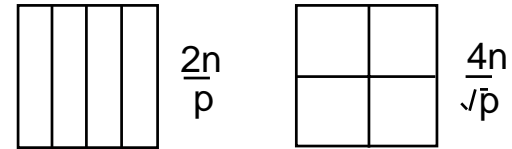
**2**

# Coverage: Levels of Optimization

**Example  
Grid Problem**

- Many ways in which an application can be suboptimal

- *Algorithmic*, e.g. assignment, blocking



- *Data structuring*, e.g. 2-d or 4-d arrays for SAS grid problem

- *Data layout, distribution and alignment*, even if properly structured

- *Orchestration*

- contention
- long versus short messages
- synchronization frequency and cost, ...

- Also, random problems with “unimportant” data structures

- Optimizing applications takes work

- Many practical applications may not be very well optimized

- May examine selected different levels to test robustness of system

**CMPE655 - Shaaban**

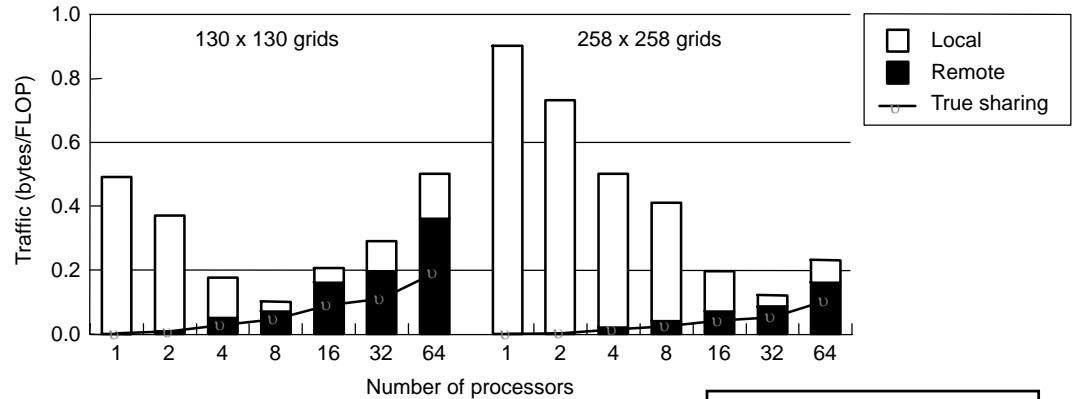
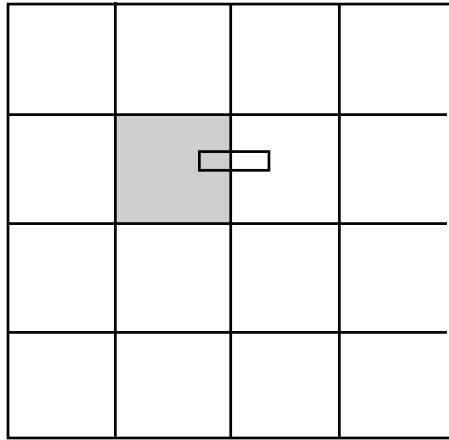


# Desirable Properties of Workloads:

## 3 Concurrency

- **Should have enough to utilize the processors**
  - If load imbalance dominates, may not be much machine can do
  - (Still, useful to know what kinds of workloads/configurations don't have enough concurrency)
- **Algorithmic speedup**: useful measure of concurrency/imbalance
  - Speedup (under scaling model) assuming all memory/communication operations take zero time
  - Ignores memory system, measures imbalance and extra work
  - Uses PRAM machine model (Parallel Random Access Machine)
    - Unrealistic, but widely used for theoretical algorithm development
- At least, should isolate performance limitations due to program characteristics that a machine cannot do much about (concurrency) from those that it can.

# Effect of Problem Size Example: Ocean



$n$ -by- $n$  grid with  $p$  processors  
(computation like grid solver)

$n/p$  is large  $\Rightarrow$

- Low communication to computation ratio
- Good spatial locality with large cache lines
- Data distribution and false sharing not problems even with 2-d array
- Working set doesn't fit in cache; high local capacity miss rate.

$n/p$  is small  $\Rightarrow$

- High communication to computation ratio
- Spatial locality may be poor; false-sharing may be a problem
- Working set fits in cache; low capacity miss rate.

*e.g. Shouldn't make conclusions about spatial locality based only on small problems, particularly if these are not very representative.*

For block decomposition:

$$\text{Communication} = \frac{4n}{\sqrt{p}}$$

$$\text{Computation} = \frac{n^2}{p}$$

$$c-to-c = \frac{4 \times \sqrt{p}}{n}$$